



Expectation Maximization

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 9th, 2007

©2005-2007 Carlos Guestrin

Gaussian Bayes Classifier

Reminder

$$P(y = i | \mathbf{x}_j) = \frac{p(\mathbf{x}_j | y = i)P(y = i)}{p(\mathbf{x}_j)}$$

$$P(y = i | \mathbf{x}_j) \propto \underbrace{\frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}_j - \mu_i)\right]}_{\text{Gaussian likelihood}} \underbrace{P(y = i)}_{\text{prior}}$$

Handwritten notes:
- μ_i : class mean
- Σ_i : class covariance


Next... back to Density Estimation

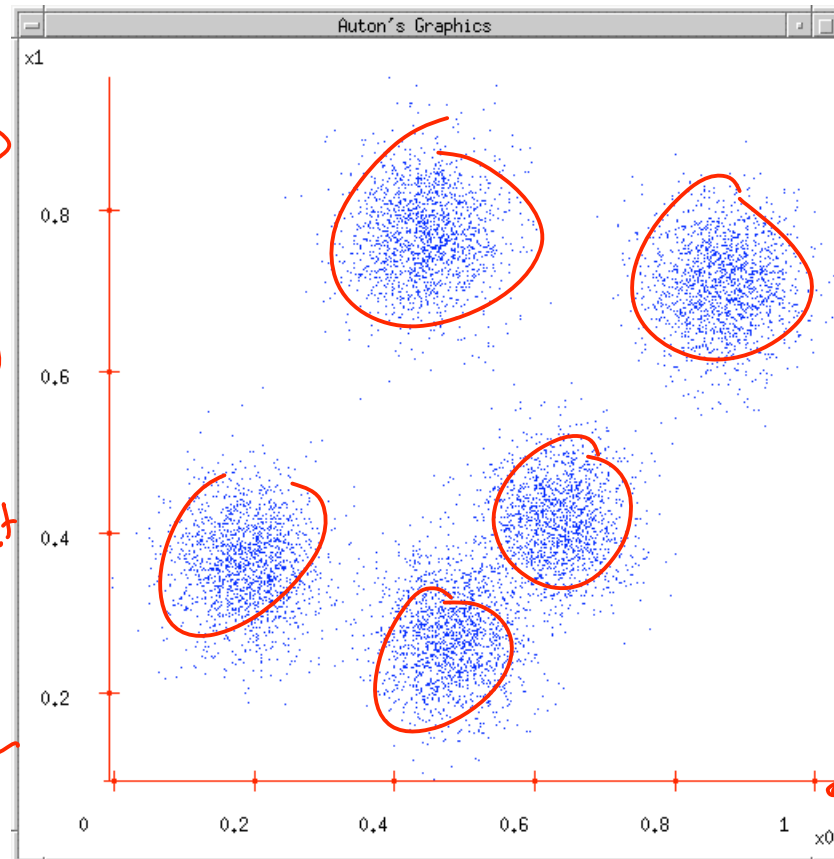
What if we want to do density estimation with multimodal or clumpy data?

want to represent $P(x)$

$$P(x) = \sum_i P(x|y=i) \cdot p(y=i)$$

↑ ↑
Gaussian weight

= 



$x = [x_0, x_1]$

Marginal likelihood for general case

$$P(y = i | \mathbf{x}_j) \propto \frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}_j - \mu_i)\right] P(y = i)$$

■ Marginal likelihood:

log

$$\prod_{j=1}^m P(\mathbf{x}_j) = \prod_{j=1}^m \sum_{i=1}^k P(\mathbf{x}_j, y = i)$$

$$= \prod_{j=1}^m \sum_{i=1}^k \frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}_j - \mu_i)\right] P(y = i)$$

$$= \sum_{j=1}^m \log \sum_{i=1}^k$$

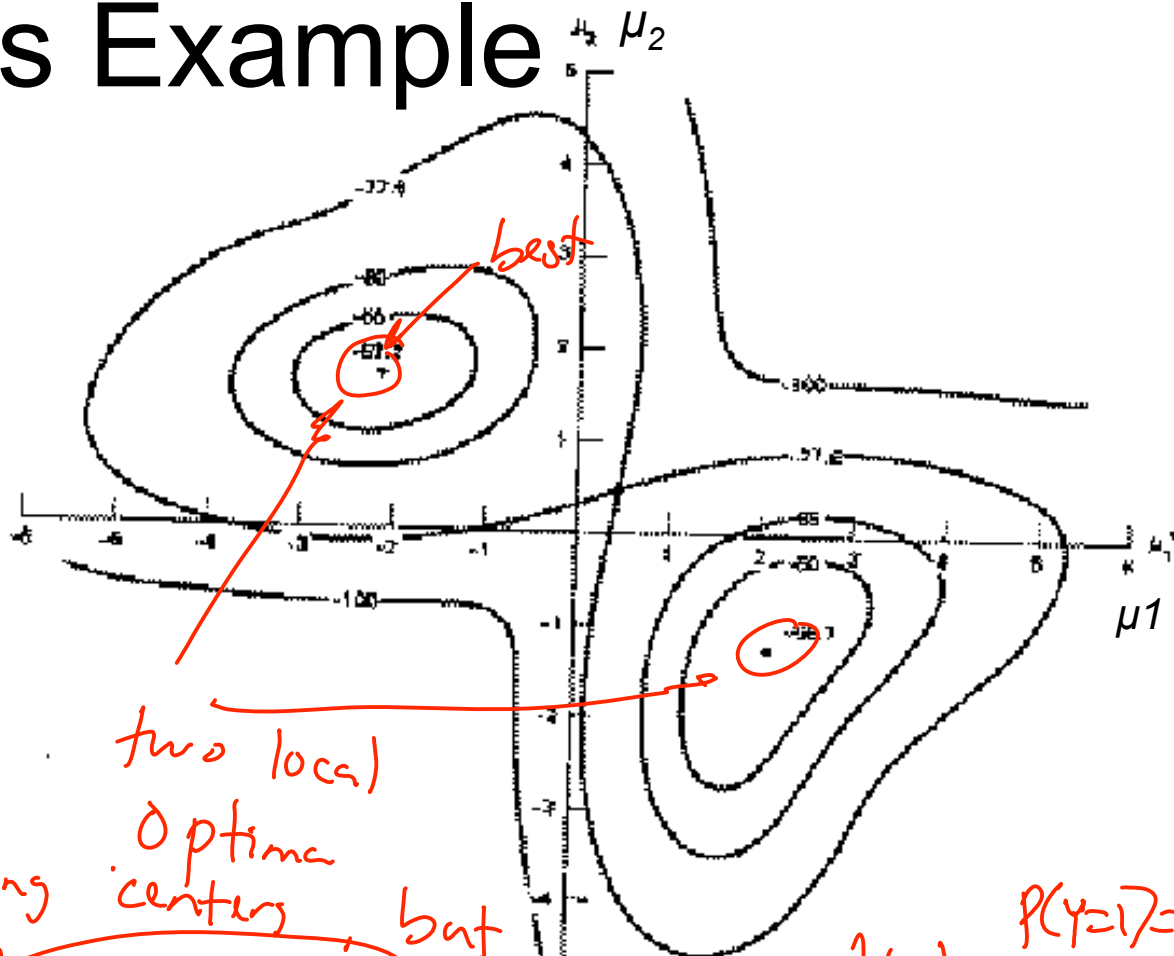
defn. of GMM $x_j \leftarrow$ observed
 assumption: $x_j \sim$ GMM
 $P(x_j) = \sum P(y=i) \cdot P(x_j | y=i)$
 don't observe y_j \Rightarrow max $P(x_j)$ Gaussian

Duda & Hart's Example



$\max_{\mu} P(X|\mu)$

Graph of
 $\log P(x_1, x_2 \dots x_{25} | \mu_1, \mu_2)$
 against $\mu_1 (\rightarrow)$ and $\mu_2 (\uparrow)$



Max likelihood = $(\mu_1 = -2.13, \mu_2 = 1.668)$

Local minimum, but very close to global at $(\mu_1 = 2.085, \mu_2 = -1.257)^*$

* corresponds to switching y_1 with y_2 .

$P(y=1) = \frac{2}{3}$
 $P(y=2) = \frac{1}{3}$

Finding the max likelihood $\mu_1, \mu_2 \dots \mu_k$



We can compute $P(\text{data} \mid \mu_1, \mu_2 \dots \mu_k)$

How do we find the μ_i 's which give max. likelihood?

- The normal max likelihood trick:

Set $\frac{\partial}{\partial \mu_i} \log \text{Prob} (\dots) = 0$

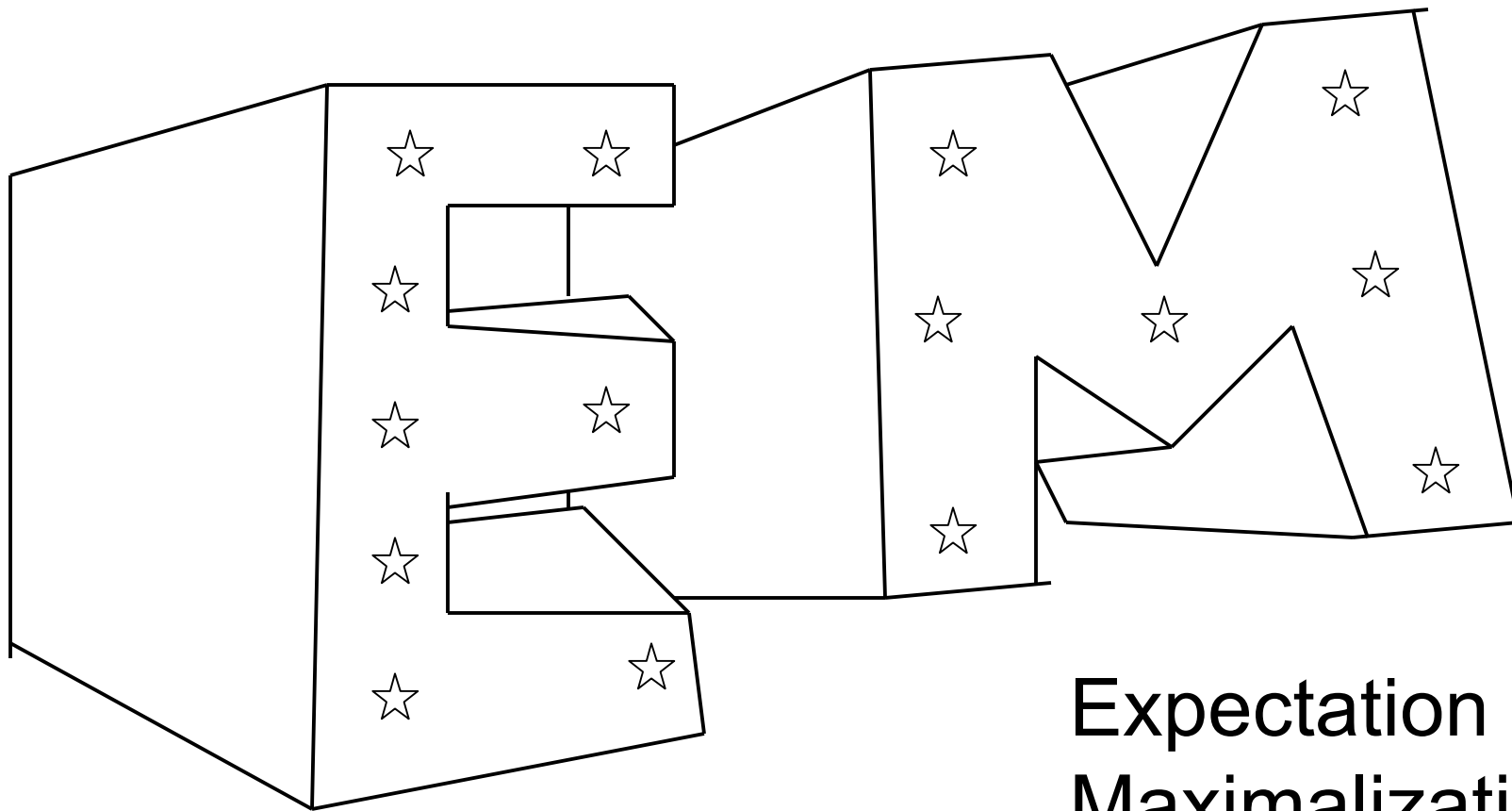
and solve for μ_i 's.

Here you get non-linear non-analytically-solvable equations

- Use gradient descent

Slow but doable

- Use a much faster, cuter, and recently very popular method...



Expectation
Maximalization

The E.M. Algorithm



DETOUR

- We'll get back to unsupervised learning soon
- But now we'll look at an even simpler case with hidden information
- The EM algorithm
 - Can do trivial things, such as the contents of the next few slides
 - An excellent way of doing our unsupervised learning problem, as we'll see
 - Many, many other uses, including learning BNs with hidden data

Silly Example

Let events be “grades in a class”

w_1 = Gets an A

$$P(A) = \frac{1}{2}$$

w_2 = Gets a B

$$P(B) = \mu$$

w_3 = Gets a C

$$P(C) = 2\mu$$

w_4 = Gets a D

$$P(D) = \frac{1}{2} - 3\mu$$

(Note $0 \leq \mu \leq 1/6$)

Assume we want to estimate μ from data. In a given class there were

a A's

b B's

c C's

d D's

What's the maximum likelihood estimate of μ given a,b,c,d ?

Trivial Statistics

$$P(A) = \frac{1}{2} \quad P(B) = \mu \quad P(C) = 2\mu \quad P(D) = \frac{1}{2} - 3\mu$$

$$P(a, b, c, d \mid \mu) = K \left(\frac{1}{2}\right)^a (\mu)^b (2\mu)^c \left(\frac{1}{2} - 3\mu\right)^d$$

$$\log P(a, b, c, d \mid \mu) = \log K + a \log \frac{1}{2} + b \log \mu + c \log 2\mu + d \log \left(\frac{1}{2} - 3\mu\right)$$

$$\text{FOR MAX LIKE } \mu, \text{ SET } \frac{\partial \text{LogP}}{\partial \mu} = 0$$

$$\frac{\partial \text{LogP}}{\partial \mu} = \frac{b}{\mu} + \frac{2c}{2\mu} - \frac{3d}{1/2 - 3\mu} = 0$$

$$\text{Gives max like } \mu = \frac{b + c}{6(b + c + d)}$$

So if class got

A	B	C	D
14	6	9	10

$$\text{Max like } \mu = \frac{1}{10}$$

Boring, but true!

Same Problem with Hidden Information

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

Someone tells us that

Number of High grades (A's + B's) = h

Number of C's = c

Number of D's = d

What is the max. like estimate of μ now?

Same Problem with Hidden Information

Someone tells us that

Number of High grades (A's + B's) = h

Number of C's = c

Number of D's = d

What is the max. like estimate of μ now?

We can answer this question circularly:

EXPECTATION

If we know the value of μ we could compute the expected value of a and b

Since the ratio $a:b$ should be the same as the ratio $\frac{1}{2} : \mu$

$$a = \frac{\frac{1}{2}}{\frac{1}{2} + \mu} h \quad b = \frac{\mu}{\frac{1}{2} + \mu} h$$

MAXIMIZATION

If we know the expected values of a and b we could compute the maximum likelihood value of μ

$$\mu = \frac{b + c}{6(b + c + d)}$$

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

E.M. for our Trivial Problem

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

We begin with a guess for μ

We iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of μ and a and b .

Define $\mu^{(t)}$ the estimate of μ on the t 'th iteration

$b^{(t)}$ the estimate of b on t 'th iteration

$\mu^{(0)}$ = initial guess

$$b^{(t)} = \frac{\mu^{(t)} h}{\frac{1}{2} + \mu^{(t)}} = E[b | \mu^{(t)}]$$

$$\mu^{(t+1)} = \frac{b^{(t)} + c}{6(b^{(t)} + c + d)}$$

= max like est. of μ given $b^{(t)}$

E-step

M-step

Continue iterating until converged.

Good news: Converging to local optimum is assured.

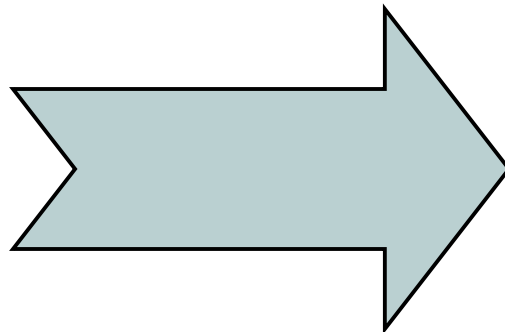
Bad news: I said "local" optimum.

E.M. Convergence

- Convergence proof based on fact that $\text{Prob}(\text{data} \mid \mu)$ must increase or remain same between each iteration [NOT OBVIOUS]
 - But it can never exceed 1 [OBVIOUS]
- So it must therefore converge [OBVIOUS]

In our example,
suppose we had

$$\begin{aligned}h &= 20 \\c &= 10 \\d &= 10 \\\mu^{(0)} &= 0\end{aligned}$$



Convergence is generally linear: error decreases by a constant factor each time step.

t	$\mu^{(t)}$	$b^{(t)}$
0	0	0
1	0.0833	2.857
2	0.0937	3.158
3	0.0947	3.185
4	0.0948	3.187
5	0.0948	3.187
6	0.0948	3.187

Back to Unsupervised Learning of GMMs – a simple case

A simple case:

We have unlabeled data $\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_m$

We know there are k classes

We know $P(y_1) P(y_2) P(y_3) \dots P(y_k)$

We don't know $\mu_1 \mu_2 \dots \mu_k$

We can write $P(\text{data} \mid \mu_1 \dots \mu_k)$

$$= p(x_1 \dots x_m \mid \mu_1 \dots \mu_k)$$

$$= \prod_{j=1}^m p(x_j \mid \mu_1 \dots \mu_k)$$

$$= \prod_{j=1}^m \sum_{i=1}^k p(x_j \mid \mu_i) P(y = i)$$

$$\propto \prod_{j=1}^m \sum_{i=1}^k \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y = i)$$

EM for simple case of GMMs: The E-step

- If we know $\mu_1, \dots, \mu_k \rightarrow$ easily compute prob. point x_j belongs to class $y=i$

$$p(y = i | x_j, \mu_1 \dots \mu_k) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y = i)$$

EM for simple case of GMMs: The M-step

- If we know prob. point x_j belongs to class $y=i$
 - MLE for μ_i is weighted average
- imagine k copies of each x_j , each with weight $P(y=i|x_j)$:

$$\mu_i = \frac{\sum_{j=1}^m P(y=i|x_j) x_j}{\sum_{j=1}^m P(y=i|x_j)}$$

E.M. for GMMs

E-step

Compute “expected” classes of all datapoints for each class

$$p(y = i | x_j, \mu_1 \dots \mu_k) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y = i)$$

*Just evaluate
a Gaussian at
 x_j*

M-step

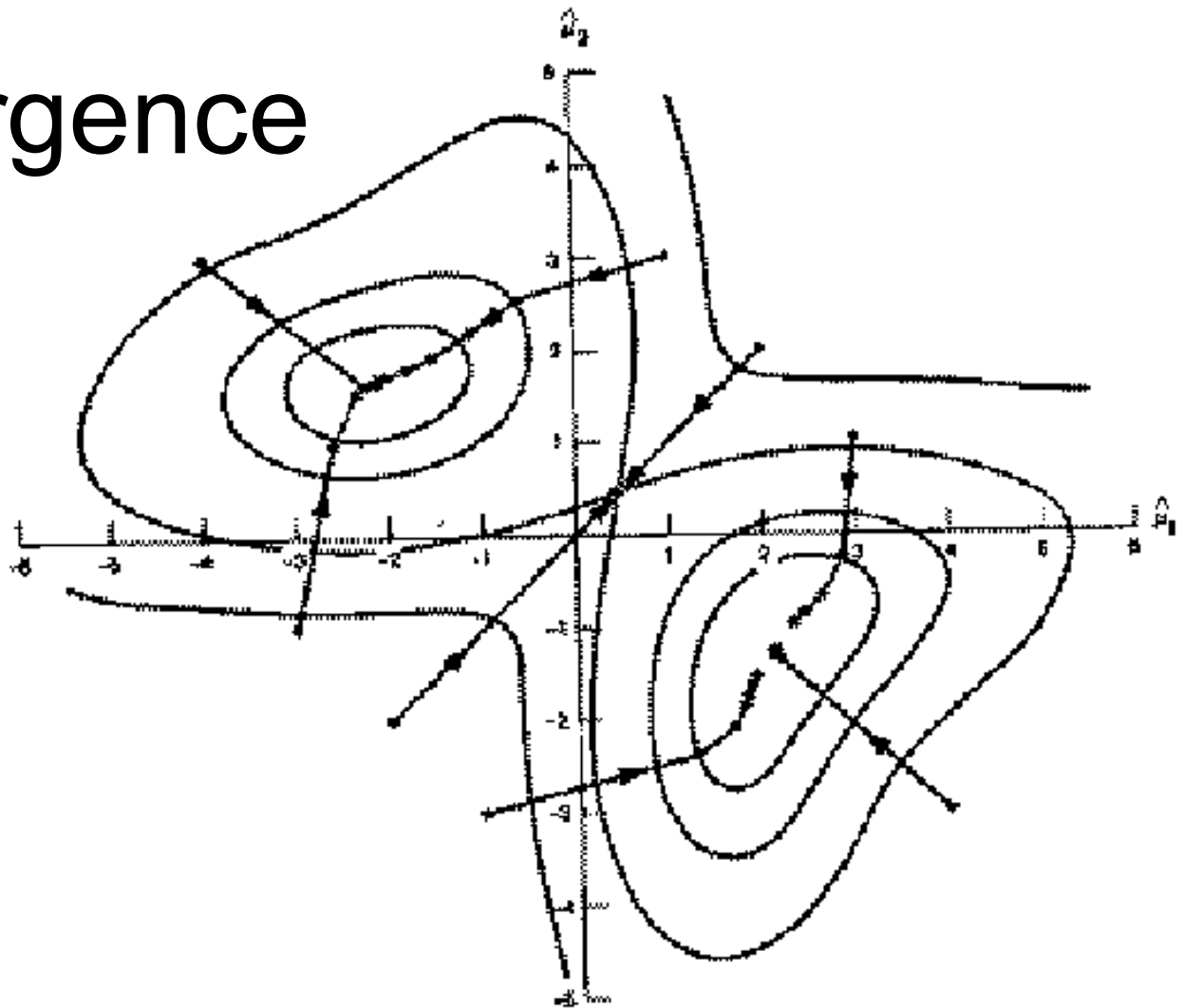
Compute Max. like μ given our data's class membership distributions

$$\mu_i = \frac{\sum_{j=1}^m P(y = i | x_j) x_j}{\sum_{j=1}^m P(y = i | x_j)}$$

E.M. Convergence



- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func. ! convergence to a local optimum guaranteed
- See Neal & Hinton reading on class webpage



- This algorithm is REALLY USED. And in high dimensional state spaces, too. E.G. Vector Quantization for Speech Data

E.M. for General GMMs

Iterate. On the t 'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \dots p_k^{(t)} \}$$

$p_i^{(t)}$ is shorthand for estimate of $P(y=i)$ on t 'th iteration

E-step

Compute “expected” classes of all datapoints for each class

$$P(y = i | x_j, \lambda_t) \propto p_i^{(t)} p(x_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

Just evaluate a Gaussian at x_j

M-step

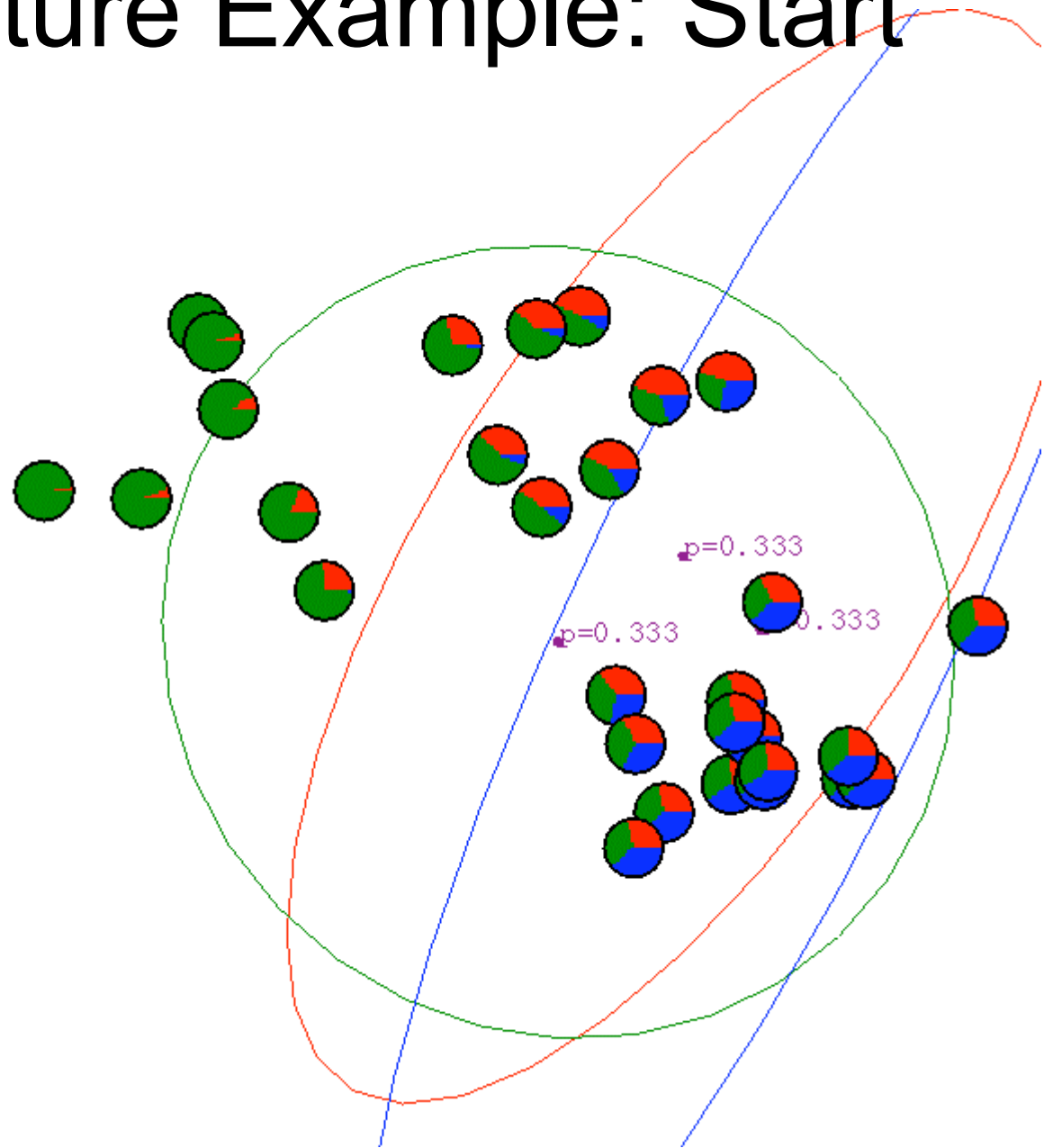
Compute Max. like μ given our data's class membership distributions

$$\hat{\mu}_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t) x_j}{\sum_j P(y = i | x_j, \lambda_t)} \quad \Sigma_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t) [x_j - \mu_i^{(t+1)}][x_j - \mu_i^{(t+1)}^T]}{\sum_j P(y = i | x_j, \lambda_t)}$$

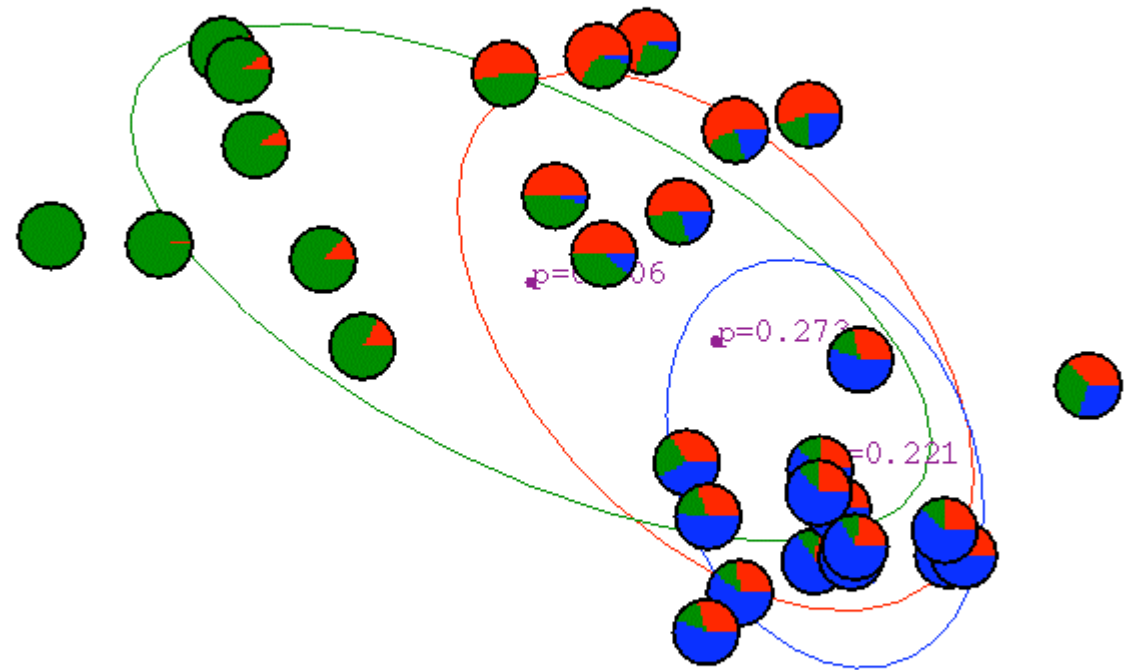
$$p_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t)}{m}$$

$m = \# \text{records}$

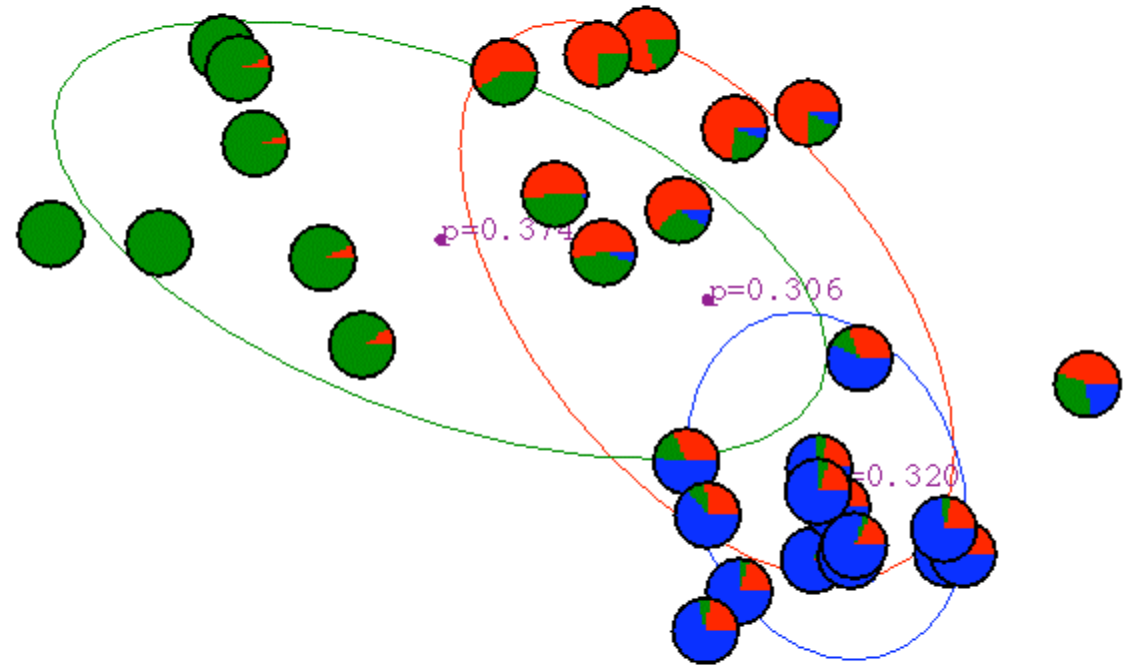
Gaussian Mixture Example: Start



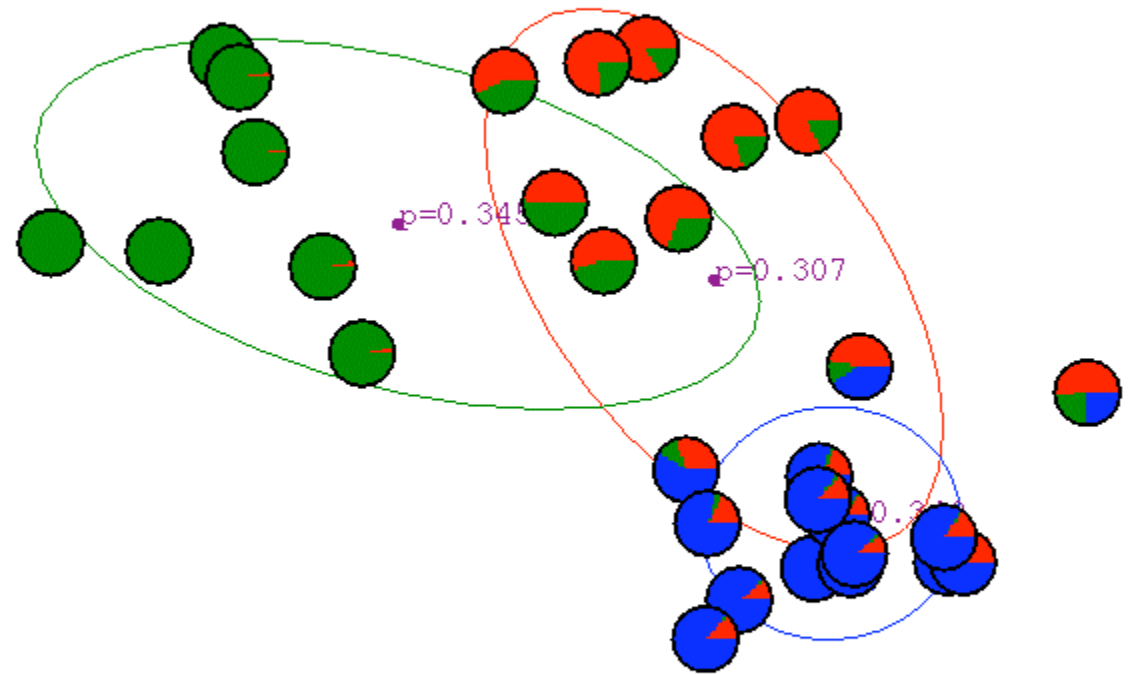
After first iteration



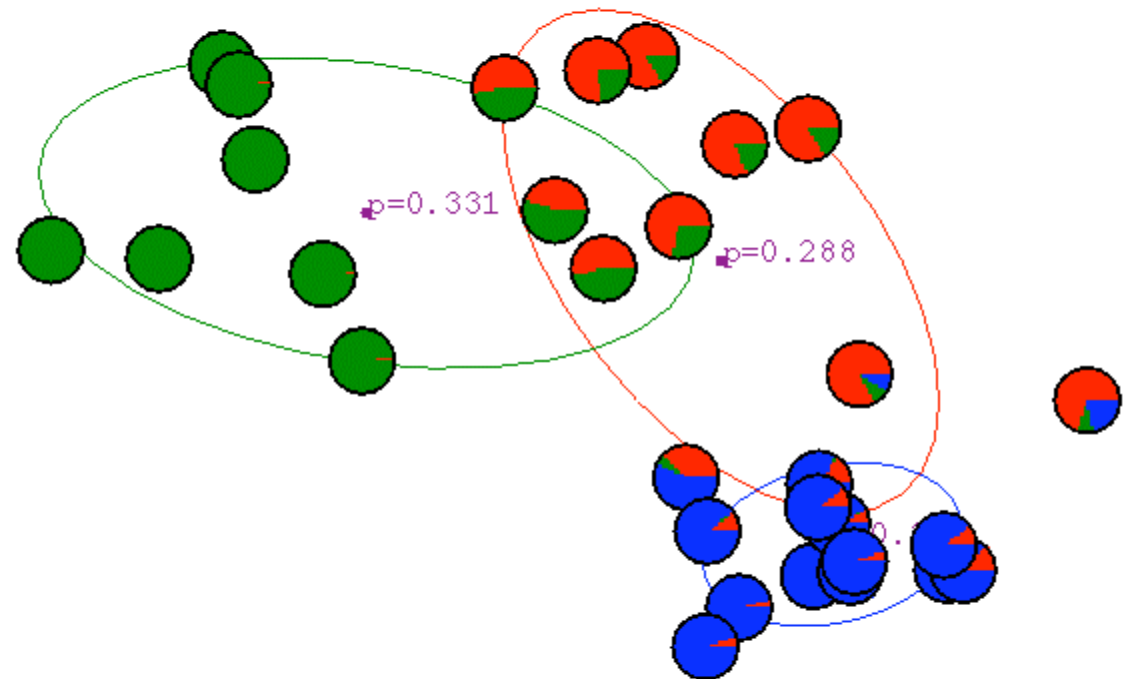
After 2nd iteration



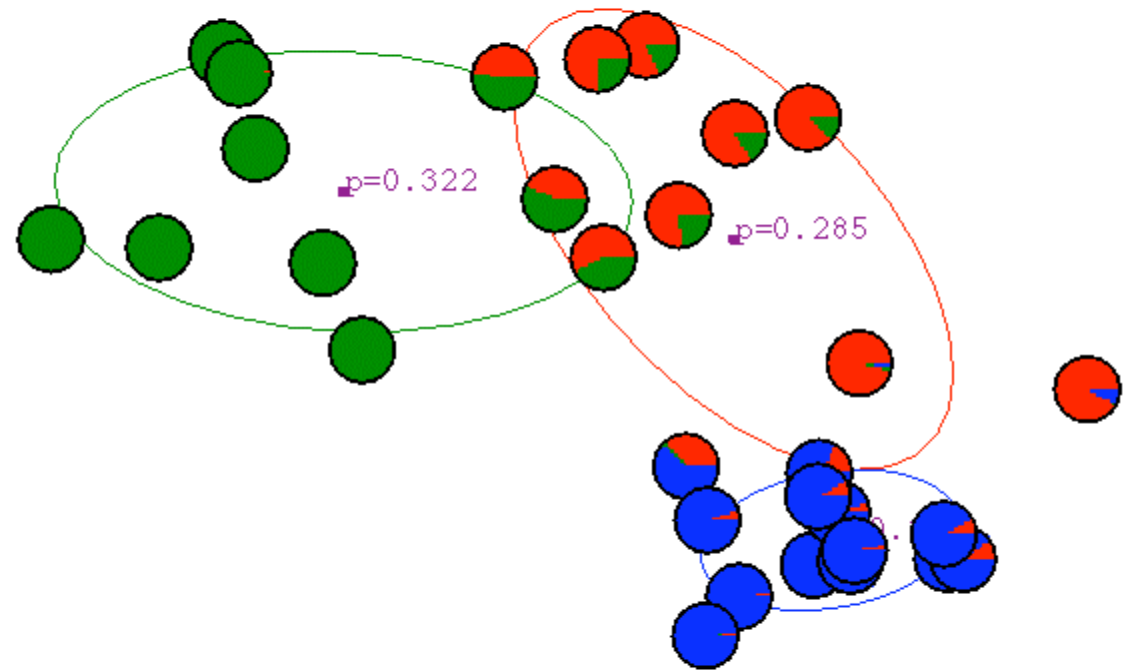
After 3rd iteration



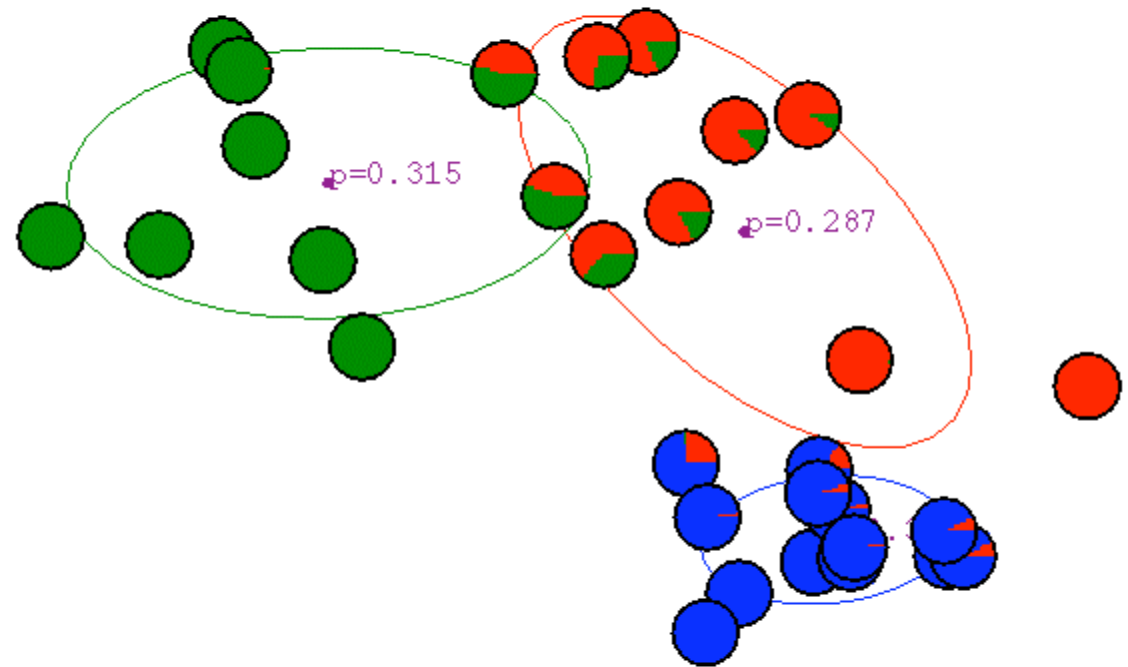
After 4th iteration



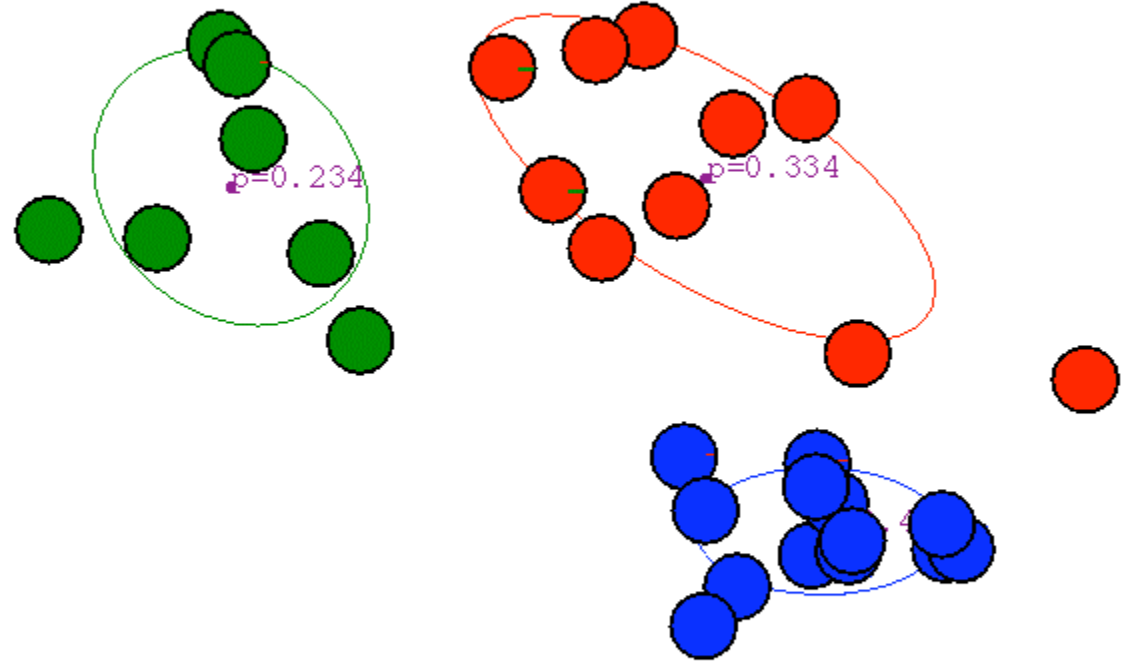
After 5th iteration



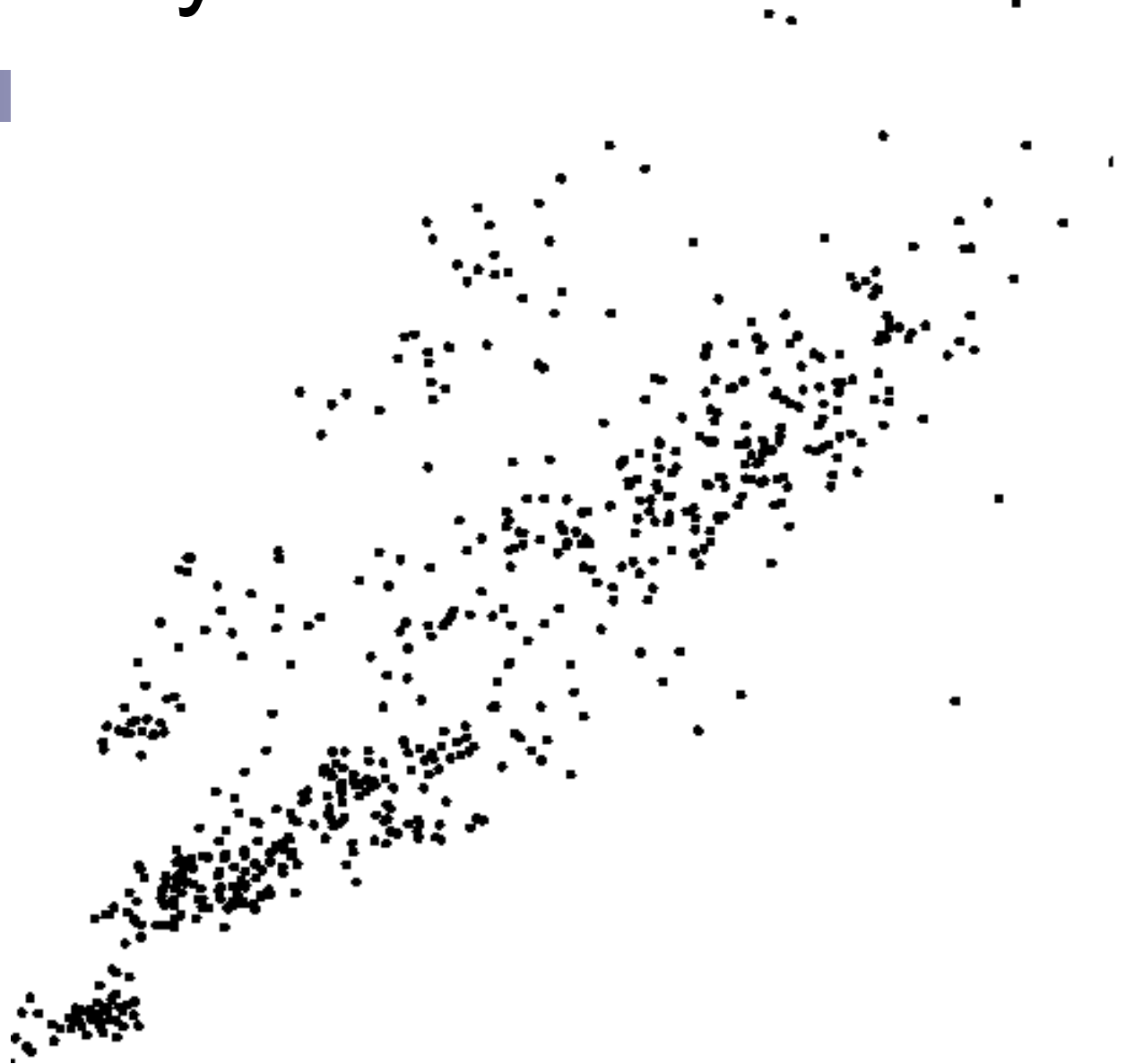
After 6th iteration



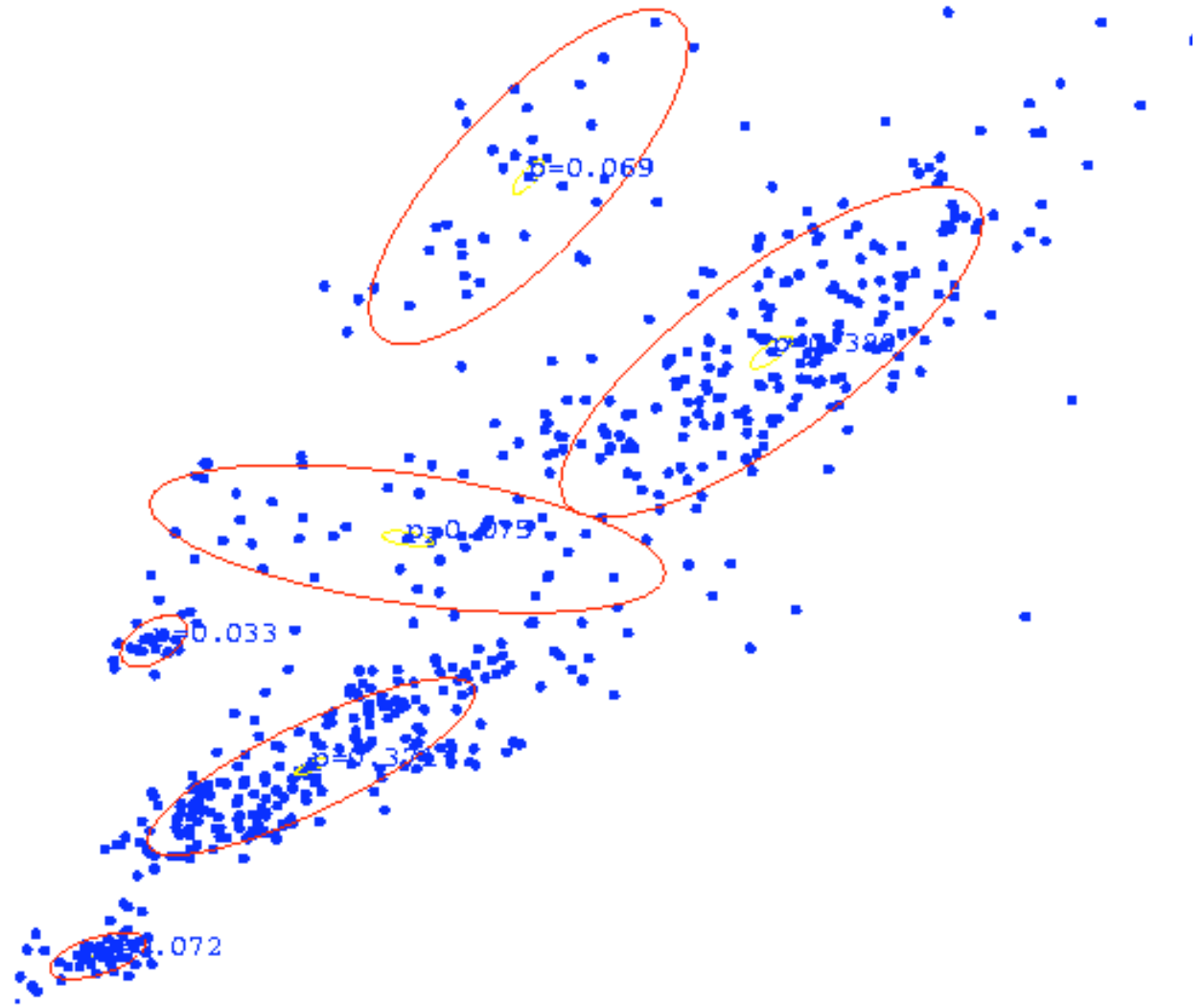
After 20th iteration

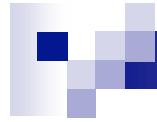


Some Bio Assay data

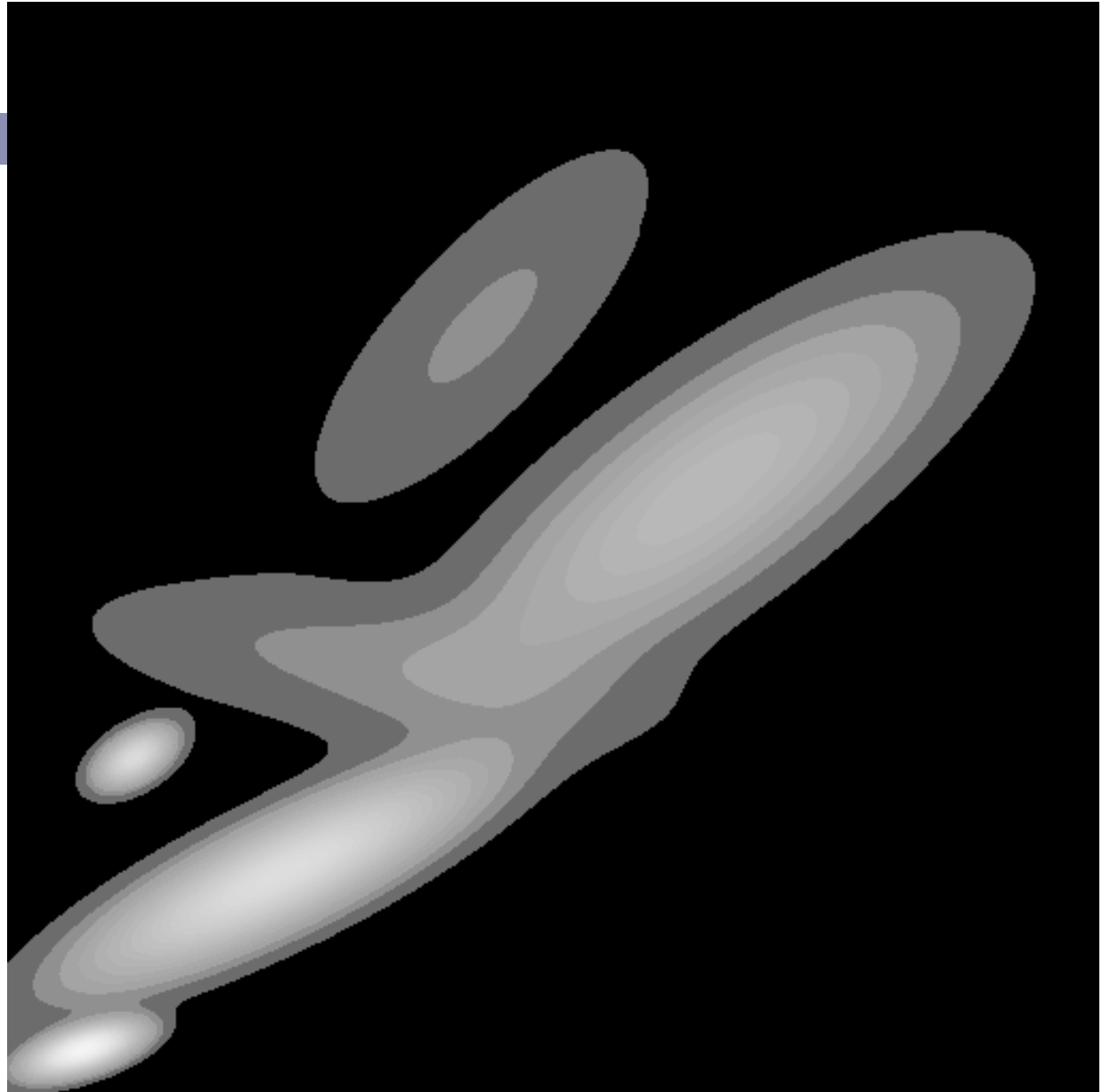


GMM clustering of the assay data





Resulting Density Estimator

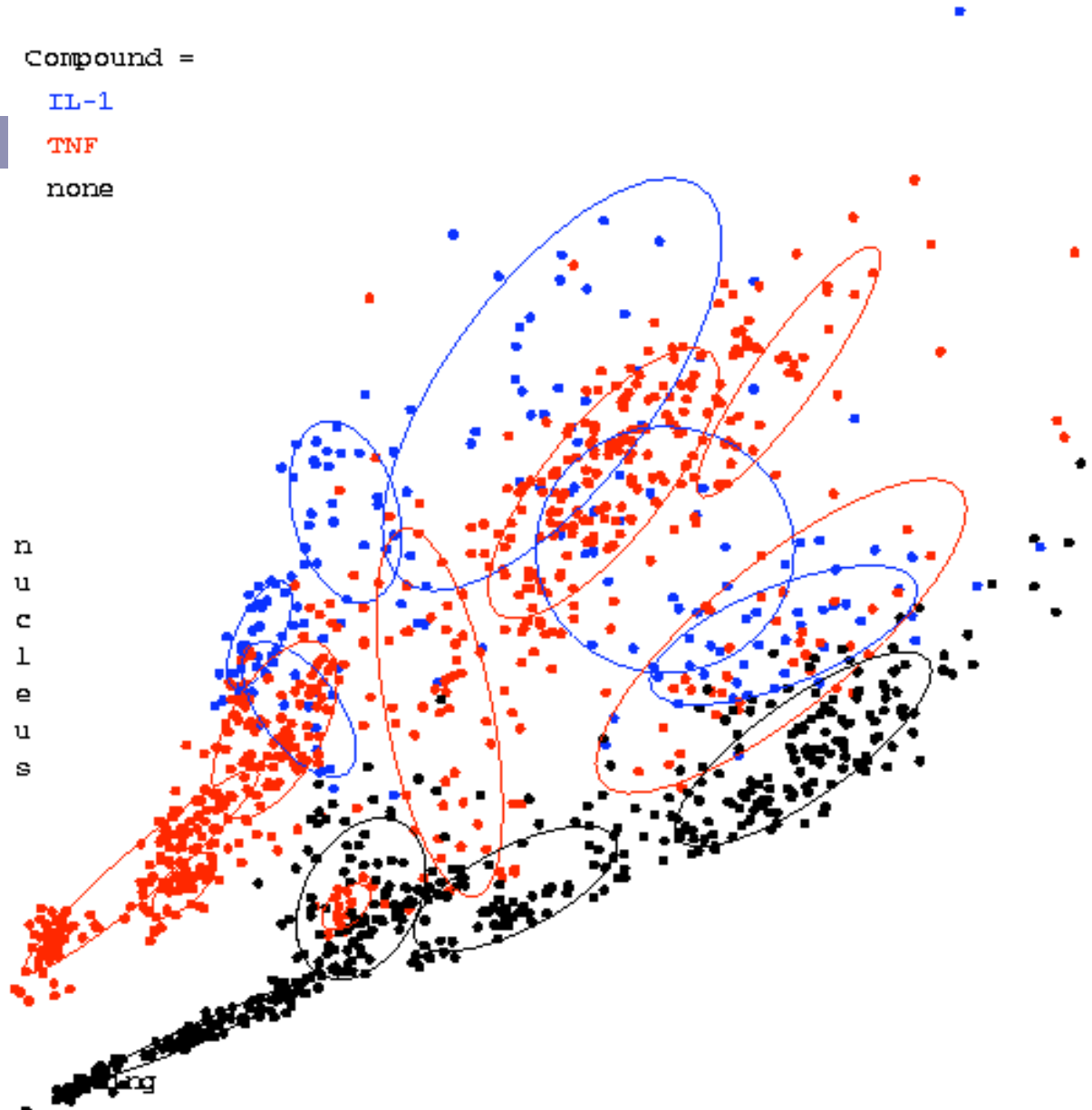




Compound =
IL-1
TNF
none

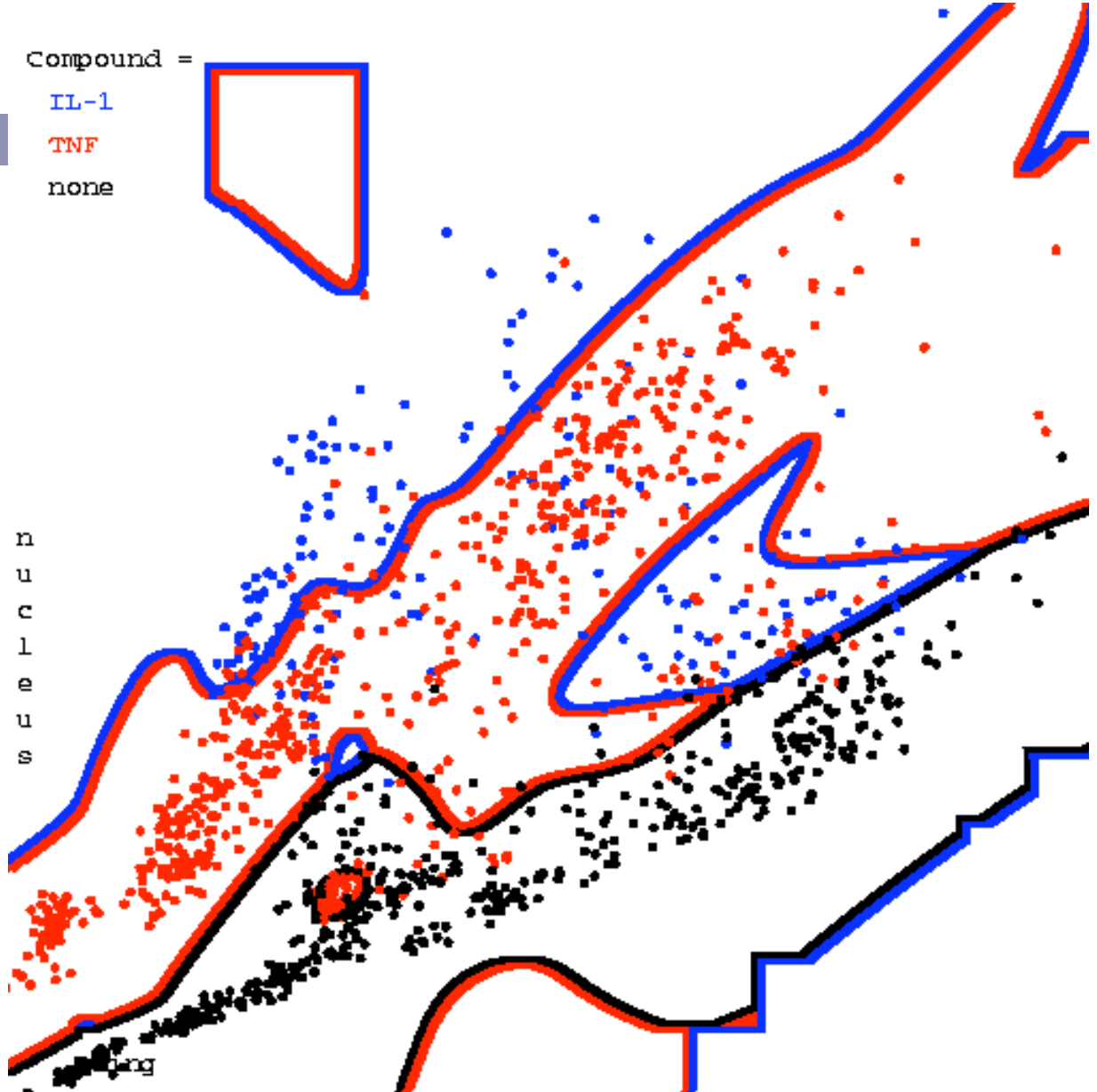
Three classes of assay


(each learned with its own mixture model)





Resulting Bayes Classifier

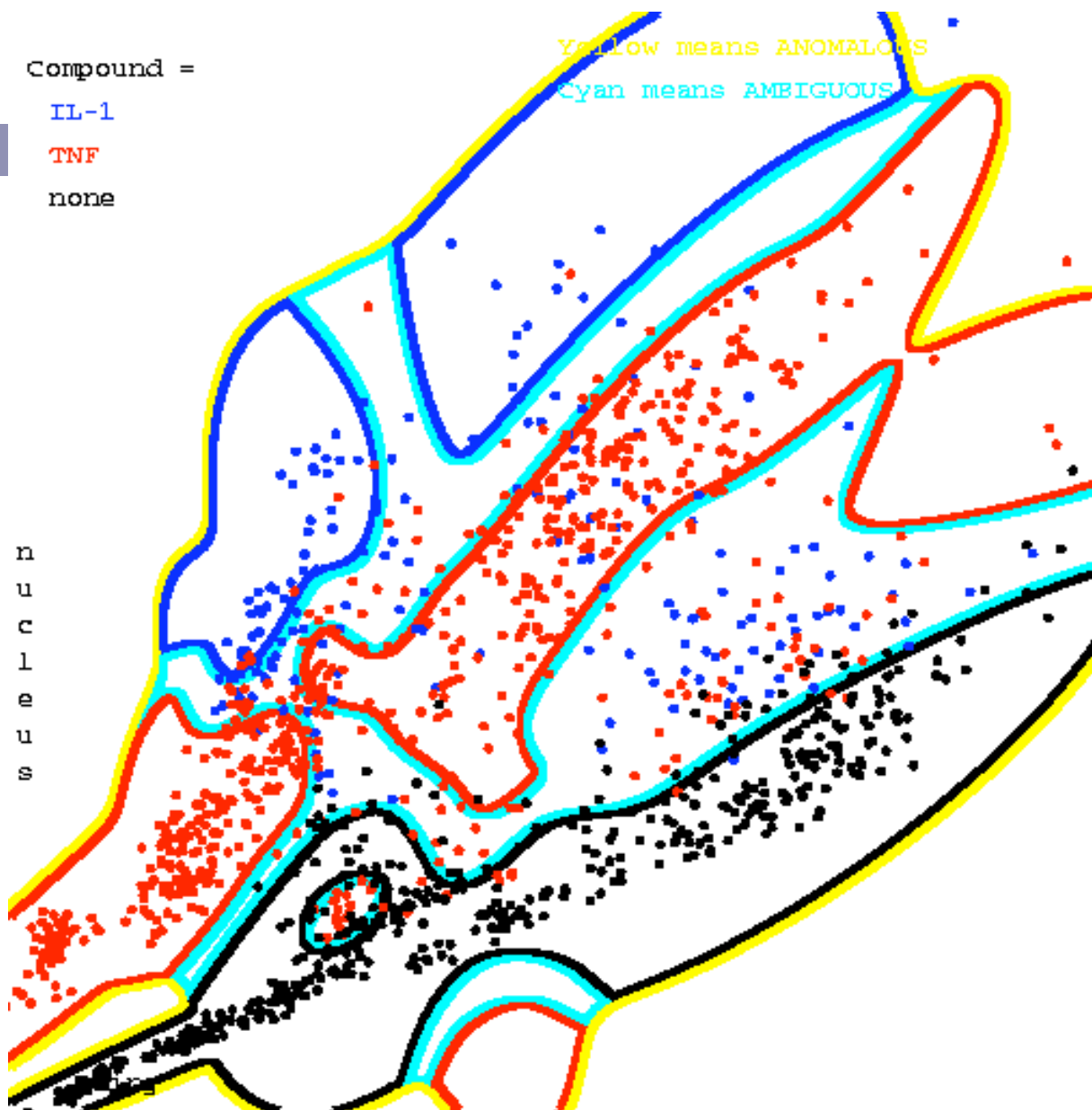




Resulting Bayes
Classifier, using
posterior
probabilities to
alert about
ambiguity and
anomalousness

**Yellow means
anomalous**

**Cyan means
ambiguous**



The general learning problem with missing data

- Marginal likelihood – \mathbf{x} is observed, \mathbf{z} is missing:


$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log \prod_{j=1}^m P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta)\end{aligned}$$

E-step

- \mathbf{x} is observed, \mathbf{z} is missing
- Compute probability of missing data given current choice of θ
 - $Q(\mathbf{z}|\mathbf{x}_j)$ for each \mathbf{x}_j
 - e.g., probability computed during classification step
 - corresponds to “classification step” in K-means

$$Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) = P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})$$

Jensen's inequality


$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{z} | \mathbf{x}_j) P(\mathbf{x}_j | \theta)$$

■ **Theorem:** $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

Applying Jensen's inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\ell(\theta^{(t)} : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)}$$

The M-step maximizes lower bound on weighted data


- Lower bound from Jensen's:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Corresponds to weighted dataset:

- ☐ $\langle \mathbf{x}_1, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 | \mathbf{x}_1)$
- ☐ $\langle \mathbf{x}_1, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2 | \mathbf{x}_1)$
- ☐ $\langle \mathbf{x}_1, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3 | \mathbf{x}_1)$
- ☐ $\langle \mathbf{x}_2, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 | \mathbf{x}_2)$
- ☐ $\langle \mathbf{x}_2, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2 | \mathbf{x}_2)$
- ☐ $\langle \mathbf{x}_2, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3 | \mathbf{x}_2)$
- ☐ ...

The M-step


$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

■ Maximization step:

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)$$

■ Use expected counts instead of counts:

- If learning requires $\text{Count}(\mathbf{x}, \mathbf{z})$
- Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{z})]$


Convergence of EM

- Define potential function $F(\theta, Q)$:

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- EM corresponds to coordinate ascent on F
 - Thus, maximizes lower bound on marginal log likelihood

M-step is easy


$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta)$$

■ Using potential function

$$F(\theta, Q^{(t+1)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t+1)})$$

E-step also doesn't decrease potential function 1

- Fixing θ to $\theta^{(t)}$:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)}$$

KL-divergence

- Measures distance between distributions

$$KL(Q||P) = \sum_z Q(z) \log \frac{Q(z)}{P(z)}$$

- KL=zero if and only if Q=P

E-step also doesn't decrease potential function 2

- Fixing θ to $\theta^{(t)}$:

$$\begin{aligned}\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) &= \ell(\theta^{(t)} : \mathcal{D}) + \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))\end{aligned}$$

E-step also doesn't decrease potential function 3

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))$$


- Fixing θ to $\theta^{(t)}$
- Maximizing $F(\theta^{(t)}, Q)$ over $Q \rightarrow$ set Q to posterior probability:

$$Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \leftarrow P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$$

- Note that

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

EM is coordinate ascent


$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- **M-step:** Fix Q , maximize F over θ (a lower bound on $\ell(\theta : \mathcal{D})$):

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m.H(Q^{(t)})$$

- **E-step:** Fix θ , maximize F over Q :

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))$$

- “Realigns” F with likelihood:

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

What you should know

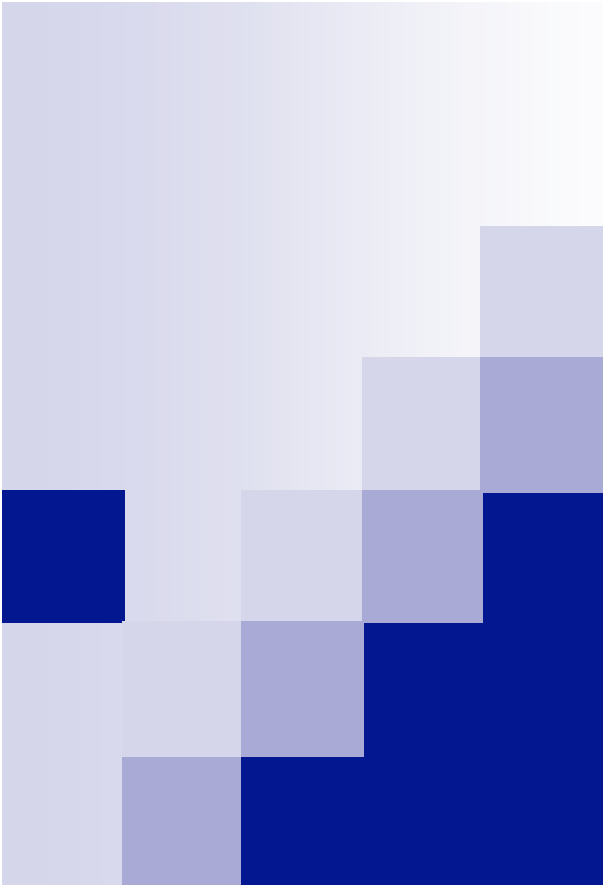


- K-means for clustering:
 - algorithm
 - converges because it's coordinate ascent
- EM for mixture of Gaussians:
 - How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent
- General case for EM

Acknowledgements



- K-means & Gaussian mixture models presentation contains material from excellent tutorial by Andrew Moore:
 - <http://www.autonlab.org/tutorials/>
- K-means Applet:
 - http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/AppletKM.html
- Gaussian mixture models Applet:
 - <http://www.neurosci.aist.go.jp/%7Eakaho/MixtureEM.html>



EM for HMMs a.k.a. The Baum-Welch Algorithm

Machine Learning – 10701/15781

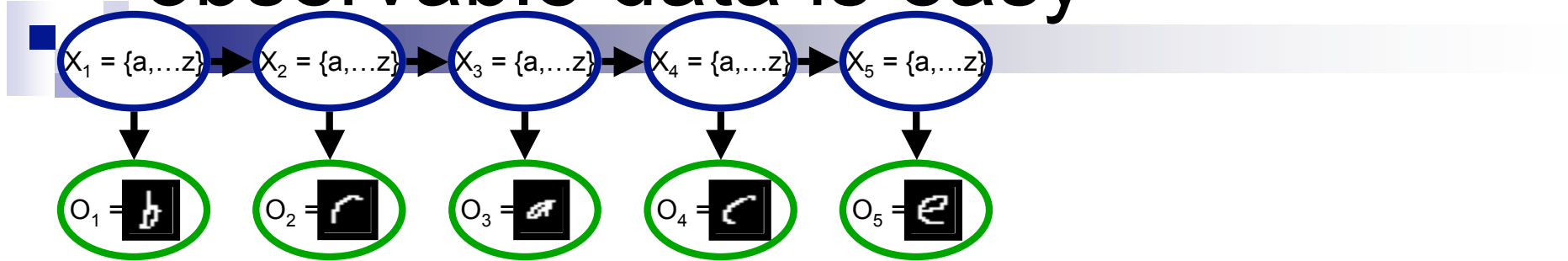
Carlos Guestrin

Carnegie Mellon University

April 9th, 2007

©2005-2007 Carlos Guestrin

Learning HMMs from fully observable data is easy



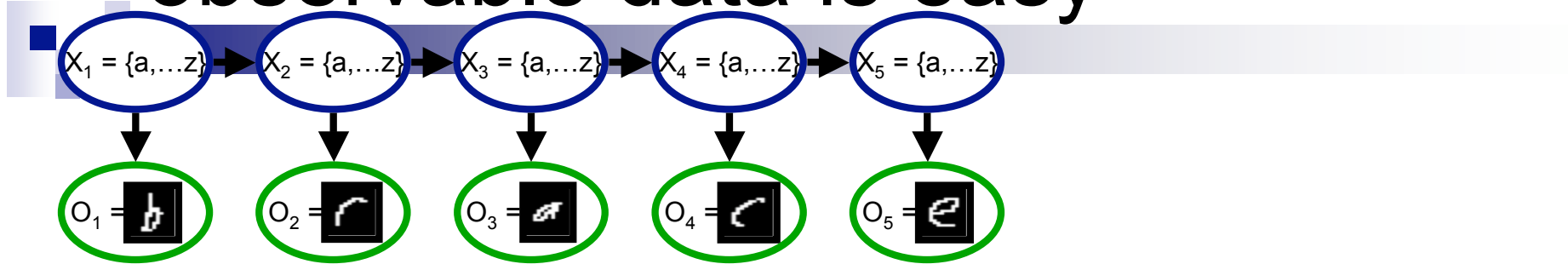
Learn 3 distributions:

$$P(X_1)$$

$$P(O_i \mid X_i)$$

$$P(X_i \mid X_{i-1})$$

Learning HMMs from fully observable data is easy



Learn 3 distributions:

$$P(X_1^a) = \frac{\text{count}(\# \text{ first letter was } a)}{N = \text{dataset size}}$$

$$P(O_i^{\text{pixel 17 is white}} | X_i^a) = \frac{\text{count}(\text{pixel 17 was white, } X_i = a)}{N_i}$$

$$P(X_i^a | X_{i-1}^b)$$

What if **O** is observed,
but **X** is hidden

select training data where letter was a

Log likelihood for HMMs when \mathbf{X} is hidden

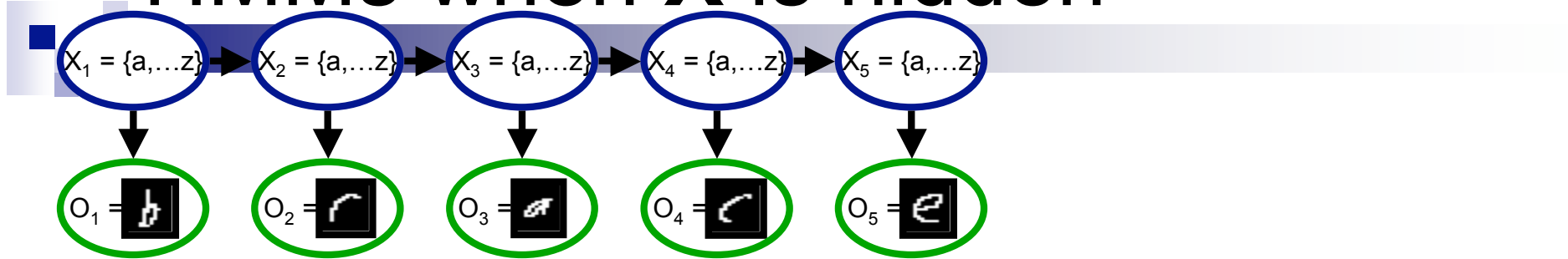
- Marginal likelihood – \mathbf{O} is observed, \mathbf{X} is missing
 - For simplicity of notation, training data consists of only one sequence:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} \mid \theta) \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} \mid \theta)\end{aligned}$$

- If there were m sequences:

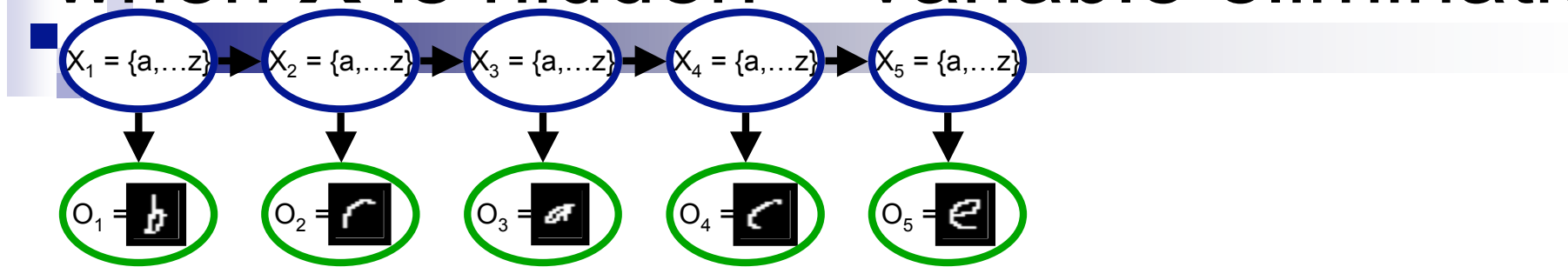
$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o}^{(j)} \mid \theta)$$

Computing Log likelihood for HMMs when **X** is hidden



$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} \mid \theta) \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} \mid \theta)\end{aligned}$$

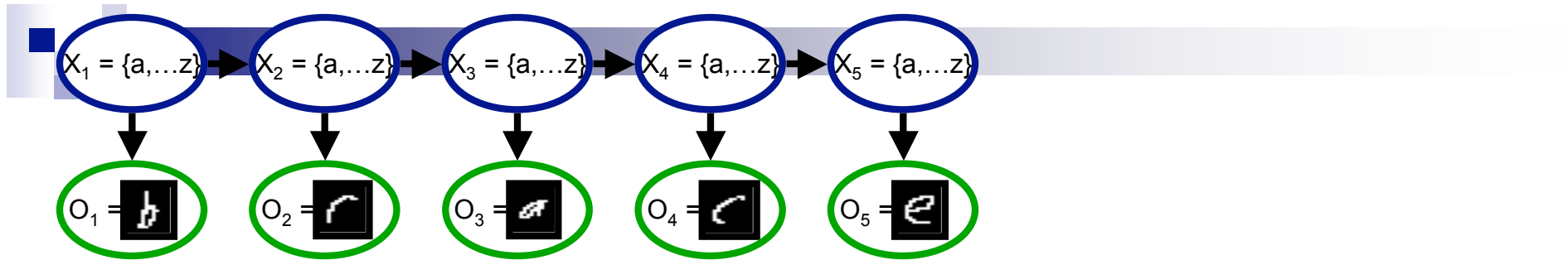
Computing Log likelihood for HMMs when **X** is hidden – variable elimination



- Can compute efficiently with variable elimination:

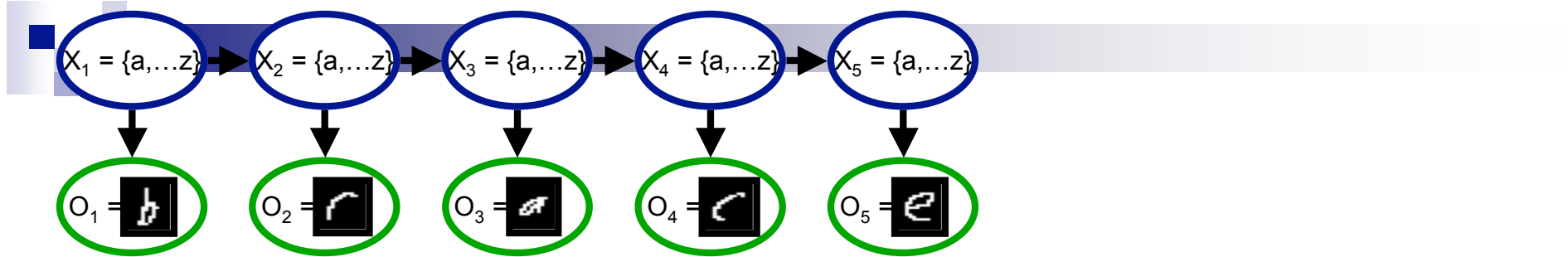
$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} \mid \theta) \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} \mid \theta)\end{aligned}$$

EM for HMMs when \mathbf{X} is hidden



- E-step: Use inference (forwards-backwards algorithm)
- M-step: Recompute parameters with weighted data

E-step

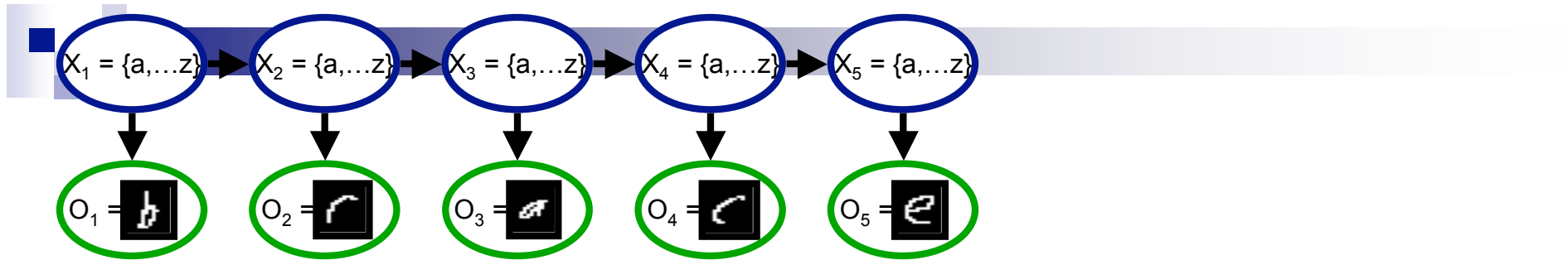


- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}

$$Q^{(t+1)}(\mathbf{x} \mid \mathbf{o}) = P(\mathbf{x} \mid \mathbf{o}, \theta^{(t)})$$

- Will correspond to inference
 - use forward-backward algorithm!

The M-step



■ Maximization step:

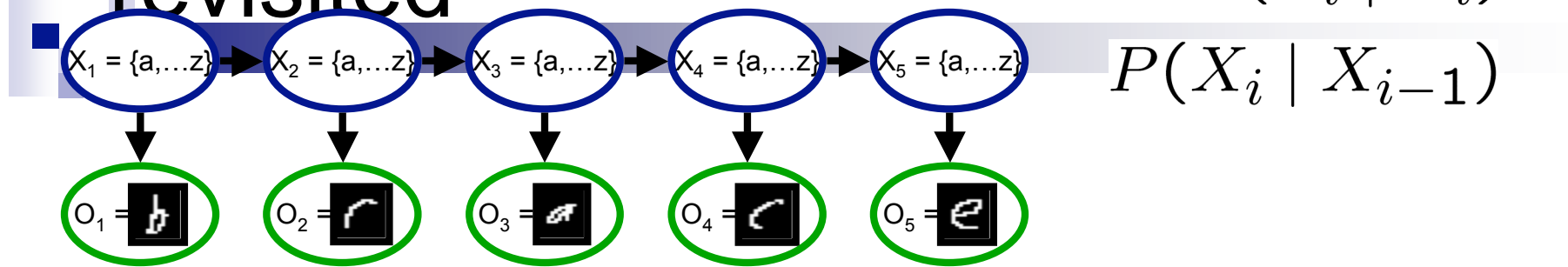
$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{\mathbf{x}} Q^{(t+1)}(\mathbf{x} \mid \mathbf{o}) \log P(\mathbf{x}, \mathbf{o} \mid \theta)$$

■ Use expected counts instead of counts:

- ☐ If learning requires $\text{Count}(\mathbf{x}, \mathbf{o})$
- ☐ Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{o})]$

Decomposition of likelihood $P(X_1)$

revisited



■ Likelihood optimization decomposes:

$$\max_{\theta} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log P(\mathbf{x}, \mathbf{o} | \theta) =$$

$$\max_{\theta} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log P(x_1 | \theta_{X_1}) P(o_1 | x_1, \theta_{O|X}) \prod_{t=2}^n P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}}) P(o_t | x_t, \theta_{O|X})$$

Starting state probability $P(X_1)$

- Using expected counts

- $P(X_1=a) = \theta_{X_1=a}$

$$\max_{\theta_{X_1}} \sum_{\mathbf{x}} Q(\mathbf{x} \mid \mathbf{o}) \log P(x_1 \mid \theta_{X_1})$$

$$\theta_{X_1=a} = \frac{\sum_{j=1}^m Q(X_1 = a \mid \mathbf{o}^{(j)})}{m}$$

Transition probability $P(X_t|X_{t-1})$

- Using expected counts

- $P(X_t=a|X_{t-1}=b) = \theta_{X_t=a|X_{t-1}=b}$

$$\max_{\theta_{X_t|X_{t-1}}} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log \prod_{t=2}^n P(x_t | x_{t-1}, \theta_{X_t|X_{t-1}})$$

$$\theta_{X_t=a|X_{t-1}=b} = \frac{\sum_{j=1}^m \sum_{t=2}^n Q(X_t = a, X_{t-1} = b | \mathbf{o}^{(j)})}{\sum_{j=1}^m \sum_{t=2}^n \sum_{i=1}^k Q(X_t = i, X_{t-1} = b | \mathbf{o}^{(j)})}$$

Observation probability $P(O_t|X_t)$

- Using expected counts

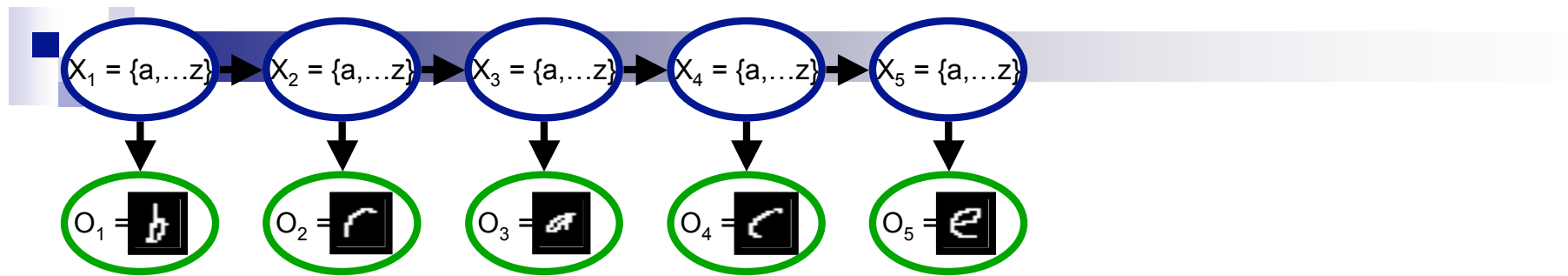
- $P(O_t=a|X_t=b) = \theta_{O_t=a|X_t=b}$

$$\max_{\theta_{O|X}} \sum_{\mathbf{x}} Q(\mathbf{x} | \mathbf{o}) \log \prod_{t=1}^n P(o_t | x_t, \theta_{O|X})$$

$$\theta_{O_t=a|X_t=b} = \frac{\sum_{j=1}^m \sum_{t=1}^n \delta(\mathbf{o}_t^{(j)} = a) Q(X_t = b | \mathbf{o}^{(j)})}{\sum_{j=1}^m \sum_{t=1}^n Q(X_t = b | \mathbf{o}^{(j)})}$$

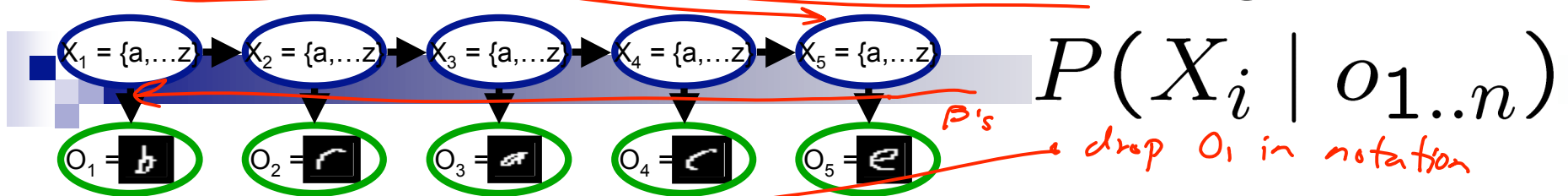
E-step revisited

$$Q^{(t+1)}(\mathbf{x} \mid \mathbf{o}) = P(\mathbf{x} \mid \mathbf{o}, \theta^{(t)})$$



- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}
- Must compute:
 - $Q(x_t = a \mid \mathbf{o})$ – marginal probability of each position
 - $Q(x_{t+1} = a, x_t = b \mid \mathbf{o})$ – joint distribution between pairs of positions

xs The forwards-backwards algorithm



■ Initialization: $\alpha_1(X_1) = P(X_1)P(o_1 | X_1)$

■ For $i = 2$ to n

□ Generate a forwards factor by eliminating X_{i-1}

sum out previous var prob obs

$$\alpha_i(X_i) = \sum_{x_{i-1}} P(o_i | X_i) P(X_i | X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

transition prob

■ Initialization: $\beta_n(X_n) = 1$

■ For $i = n-1$ to 1

□ Generate a backwards factor by eliminating X_{i+1}

xs

$$\beta_i(X_i) = \sum_{x_{i+1}} P(o_{i+1} | x_{i+1}) P(x_{i+1} | X_i) \beta_{i+1}(x_{i+1})$$

xs

■ 8 i, probability is: $P(X_i | o_{1..n}) = \alpha_i(X_i) \beta_i(X_i)$

normalized

$$\alpha_n(X_n) = P(X_n | o_{1:n})$$

normalized

$$\beta_1(X_1) \alpha_1(X_1) = P(X_1 | o_{1:n})$$

xs

$$\alpha_5(a)$$

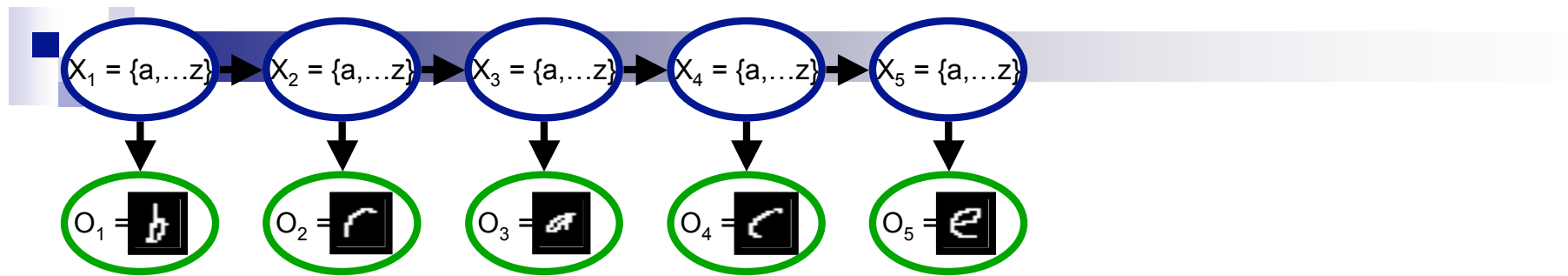
$$\alpha_5(b)$$

$$\vdots$$

$$\alpha_5(z)$$

E-step revisited

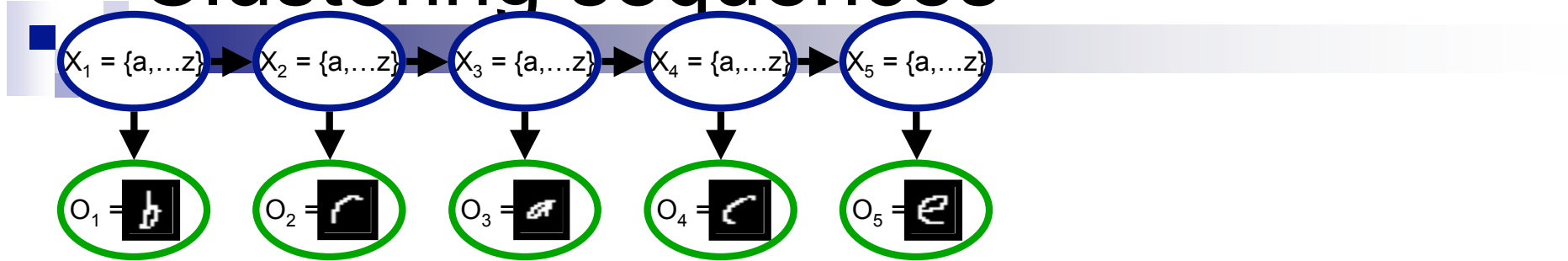
$$Q^{(t+1)}(\mathbf{x} \mid \mathbf{o}) = P(\mathbf{x} \mid \mathbf{o}, \theta^{(t)})$$



- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}
- Must compute:
 - $Q(x_t = a \mid \mathbf{o})$ – marginal probability of each position
 - Just forwards-backwards!
 - $Q(x_{t+1} = a, x_t = b \mid \mathbf{o})$ – joint distribution between pairs of positions
 - Homework! 😊

What can you do with EM for HMMs? 1

– Clustering sequences

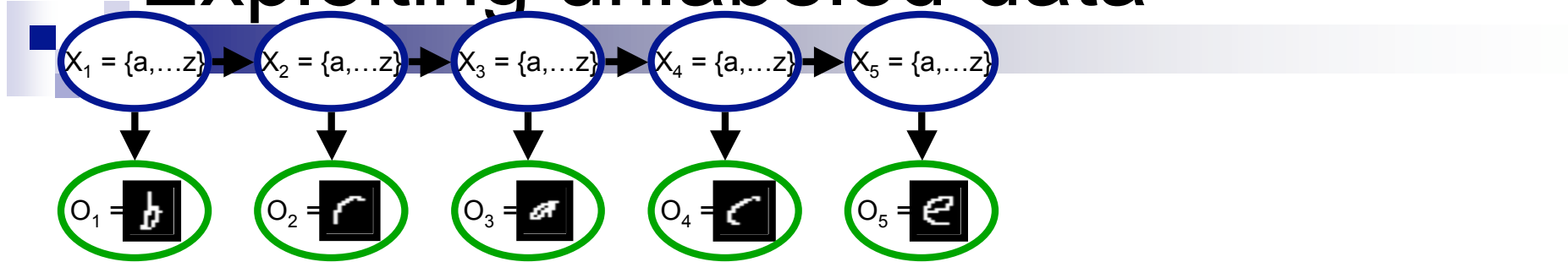


Independent clustering:

Sequence clustering:

What can you do with EM for HMMs? 2

– Exploiting unlabeled data



- Labeling data is hard work ! save (graduate student) time by using both labeled and unlabeled data

- Labeled data:

- $\langle X = \text{"brace"}, O = \text{[image of 'b']} \rangle$

- Unlabeled data:

- $\langle X = \text{?????}, O = \text{[image of 'b']} \rangle$

Exploiting unlabeled data in clustering

- A few data points are labeled

- $\langle x, o \rangle$

- Most points are unlabeled

- $\langle ?, o \rangle$

- In the E-step of EM:

- If i'th point is unlabeled:

- compute $Q(X|o_i)$ as usual

- If i'th point is labeled:

- set $Q(X=x|o_i)=1$ and $Q(X \neq x|o_i)=0$

- M-step as usual

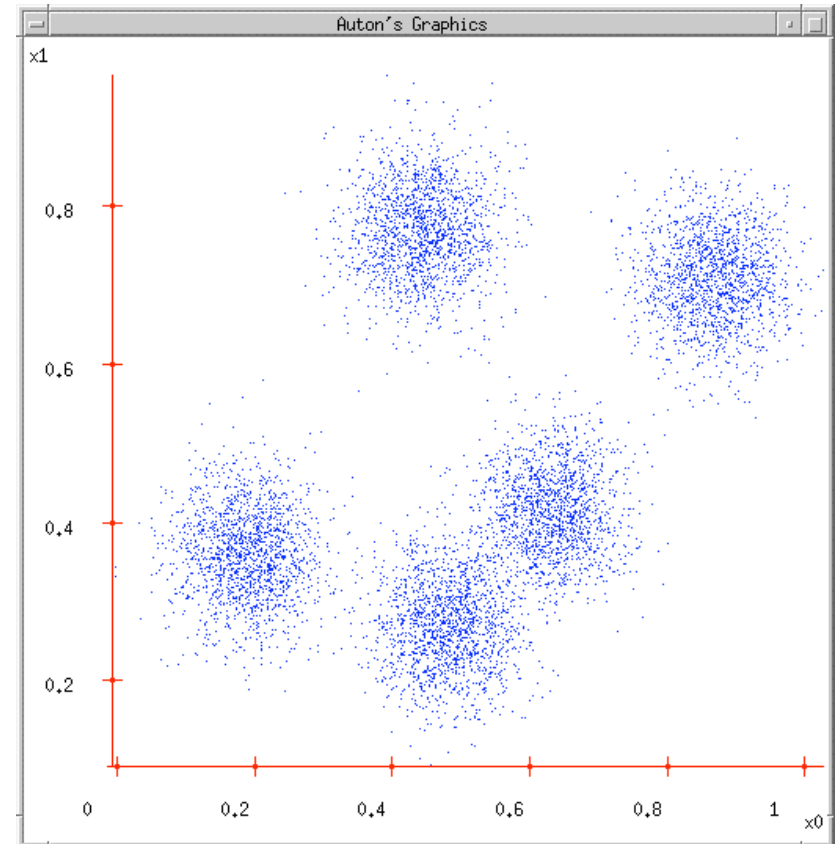
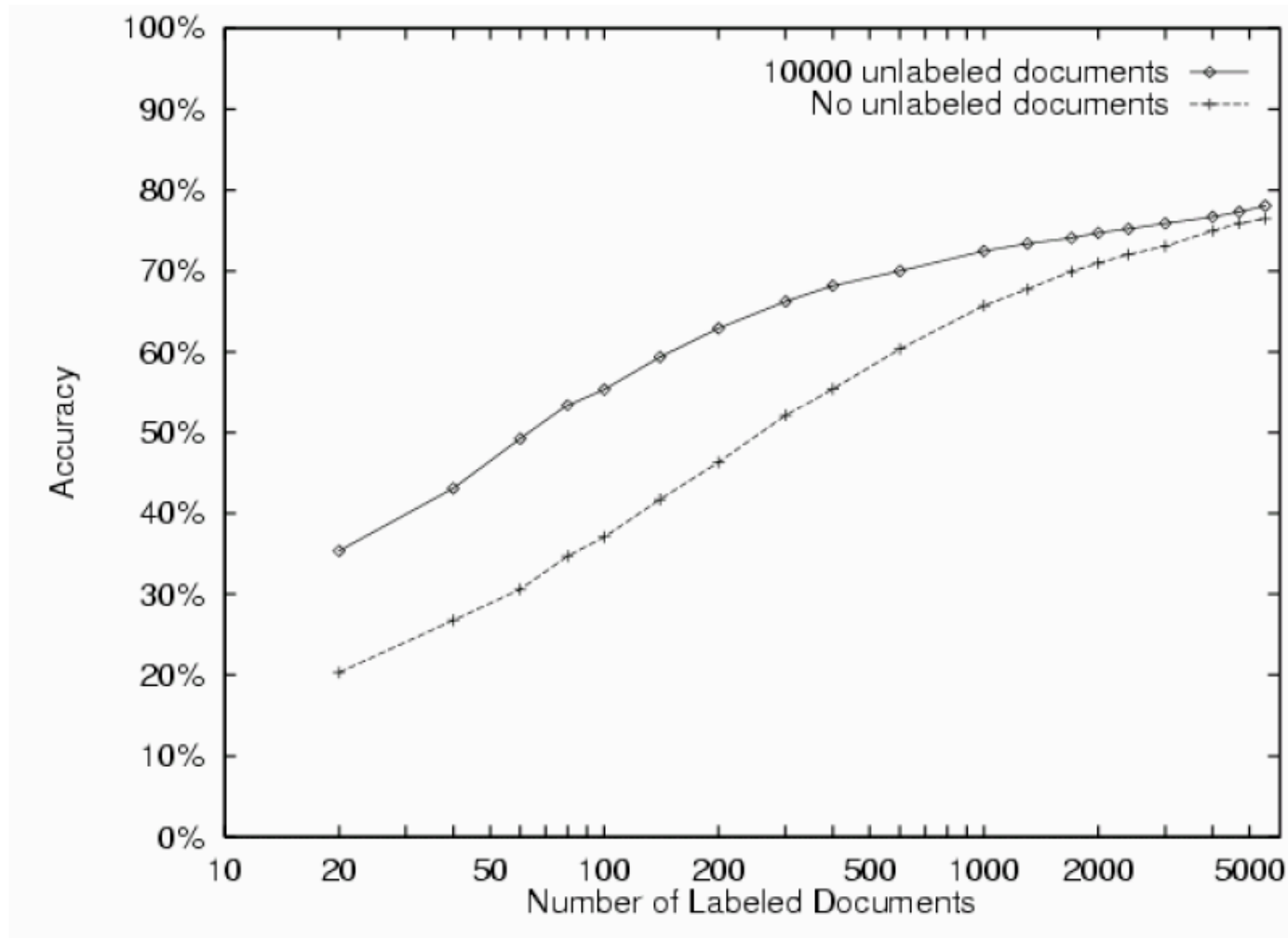


Table 3. Lists of the words most predictive of the **course** class in the WebKB data set, as they change over iterations of EM for a specific trial. By the second iteration of EM, many common **course**-related words appear. The symbol *D* indicates an arbitrary digit.

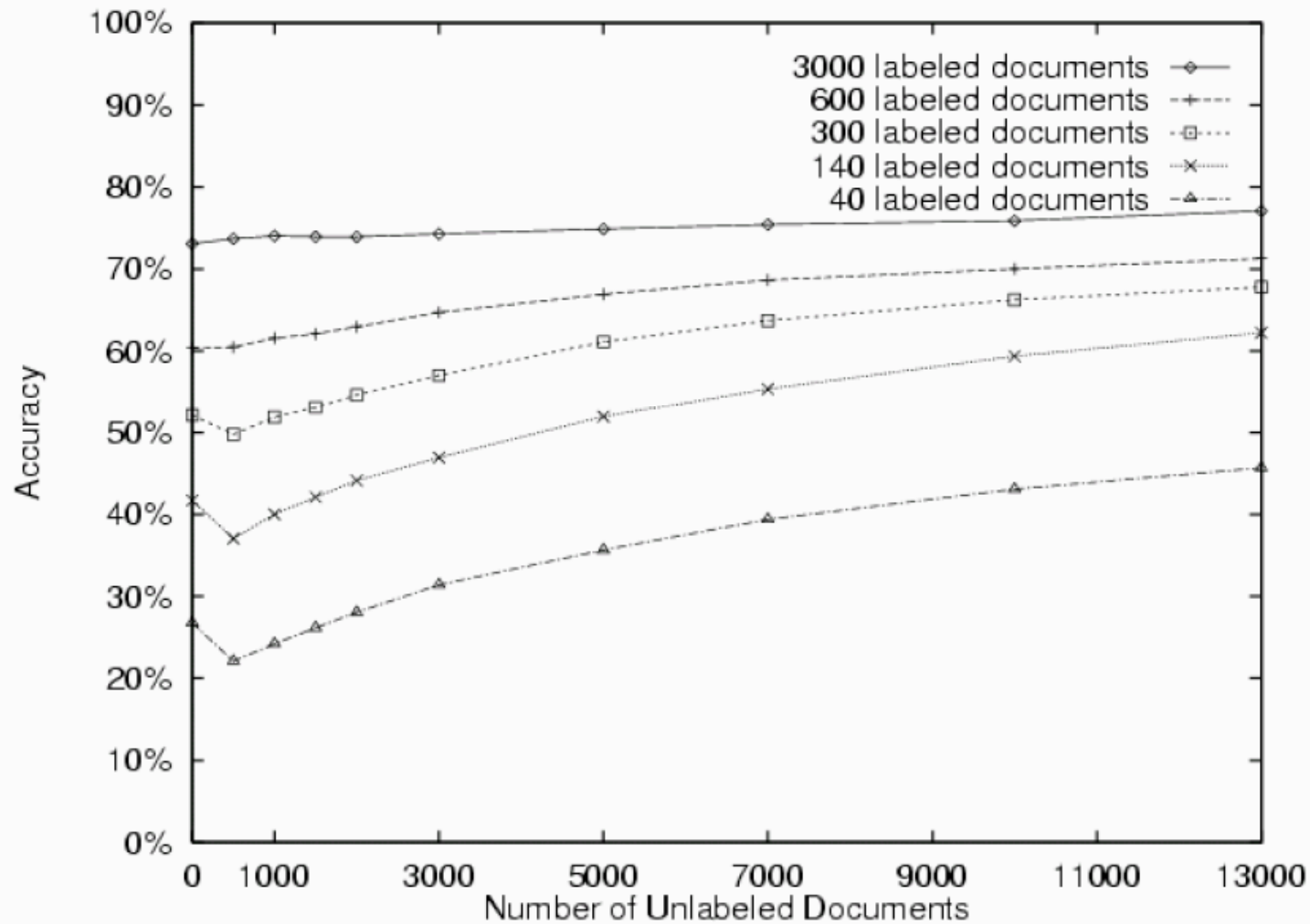
Iteration 0	Iteration 1	Iteration 2
intelligence	<i>DD</i>	<i>D</i>
<i>DD</i>	<i>D</i>	<i>DD</i>
artificial	lecture	lecture
understanding	cc	cc
<i>DDw</i>	<i>D*</i>	<i>DD:DD</i>
dist	<i>DD:DD</i>	due
identical	handout	<i>D*</i>
rus	due	homework
arrange	problem	assignment
games	set	handout
dartmouth	tay	set
natural	<i>DDam</i>	hw
cognitive	yurttas	exam
logic	homework	problem
proving	kfoury	<i>DDam</i>
prolog	sec	postscript
knowledge	postscript	solution
human	exam	quiz
representation	solution	chapter
field	assaf	ascii

Using one
labeled
example per
class

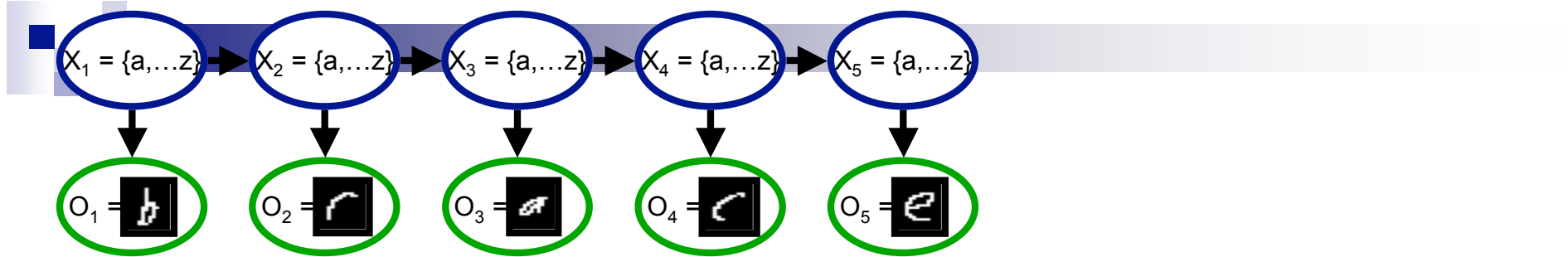
20 Newsgroups data – advantage of adding unlabeled data



20 Newsgroups data – Effect of additional unlabeled data



Exploiting unlabeled data in HMMs



- A few data points are labeled
 - $\langle x, o \rangle$
- Most points are unlabeled
 - $\langle ?, o \rangle$
- In the E-step of EM:
 - If i'th point is unlabeled:
 - compute $Q(X|o_i)$ as usual
 - If i'th point is labeled:
 - set $Q(X=x|o_i)=1$ and $Q(X \neq x|o_i)=0$
- M-step as usual
 - Speed up by remembering counts for labeled data

What you need to know



- Baum-Welch = EM for HMMs
- E-step:
 - Inference using forwards-backwards
- M-step:
 - Use weighted counts
- Exploiting unlabeled data:
 - Some unlabeled data can help classification
 - Small change to EM algorithm
 - In E-step, only use inference for unlabeled data

Acknowledgements



- Experiments combining labeled and unlabeled data provided by Tom Mitchell